

Amended agglomerative clustering for web users navigation behavior

Sumaiya Banu S¹, Kavitha K¹, Swetha V², Sathiya Raj R³

1. Department of CSE, Government College of Engineering, Bargur-635104, India
 2. Department of CSE, Kingston Engineering College, Vellore-632059, India
 3. Department of CSE, Madanapalle Institute of Technology & Science, Madanapalle, India

Received 23 March; accepted 03 May; published online 01 June; printed 16 June 2013

ABSTRACT

Clustering is a data mining technique to group a set of unsupervised data based on the conceptual clustering principal: maximizing the intra class similarity and minimizing the interclass similarity, and also the task of finding natural partitioning within a data set such that data items within the same group are more similar than those within different groups. Nowadays all available clustering techniques for web usage mining are based upon usage patterns which are derived from user's page preferences. The efficiency of usage pattern can be improved by using the similarities in user's access behavior based on time locality of user's navigational acts. Cluster shows similar visiting behavior at the same time period defined by tuning algorithm. The priority given to the time visit can be varied by user's page references.

Keywords: Web mining, Web users clustering, Navigation, Access time.

1. INTRODUCTION

Interest in the analysis of user behavior on the Web has been increasing rapidly. This increase stems from the realization that added value for Web site visitors are not gained merely through larger quantities of data on a site, but through easier access to the required information at the right time and in the most suitable form. Web Mining is the extraction of interesting and potentially useful patterns and implicit information from artifacts or activity related to the World Wide Web. Web usage mining provides the support for the web site design, providing personalization server and other business making decision, etc. In order to better serve for the users, web mining applies the data mining, the artificial intelligence and the chart technology and so on to the web data and traces users' visiting characteristics, and then extracts the users' using pattern (Qingtian Han et al. 2008). According to the differences of the mining objects, there are roughly three knowledge discovery domains that pertain to web mining: Web Content Mining, Web Structure Mining, and Web Usage Mining. Web content mining is the process of extracting knowledge from the content of documents or their descriptions. Web document text mining, resource discovery based on concepts indexing or agent; based technology may also fall in this category. Web structure mining is the process of inferring knowledge from the World Wide Web organization and links between references and referents in the Web. Finally, web usage mining, also known as Web Log Mining, is the process of extracting interesting patterns in web access logs.

1.1. Web Usage Mining

Web usage mining, also known as Web Log Mining, is the process of extracting interesting patterns in web access logs. Web servers record and accumulate data about user interactions whenever requests for resources are received. Analyzing the web access logs of different web sites can help understand the user behavior and the web structure, thereby improving the website design. Log record has lots of useful information such as URL, IP address, time and so on. Analyzing and discovering log could help us to find more potential users of the web site and trace service quality of the site. The large majority of methods that have been used for pattern discovery from Web data are clustering methods. Clustering has been used for grouping users with common browsing behavior. Web users clustering (Yan et al. 1996) is to use web access log files to partition a set of users into clusters such that the users within a cluster are more similar to each other than users from different clusters. The discovered clusters can then help in on-the-fly transformation of the web site content. In particular, web page scan be automatically linked by artificial hyperlinks. The idea is to try to match an active user's access pattern with one or more of the clusters discovered from the web log files (Rajni Pamnani et al). Pages in the matched clusters that have not been explored by the user may serve as navigational hints for the user to follow the aspects based on a weight factor. This paper we highlight the fact that grouping Web users based on their navigational.

2. EXISTING SYSTEM

More specifically, we define three different user visiting structures in order to capture all aspects of interrelations in page and time visiting. A vector is used to represent the frequency of a user's visits to particular pages (with no information about the time of visits) while a second one records the frequency of the user's visits at particular timeframes (with no information about which page was visited). Moreover, the lack of the complementary information in each of these vectors motivated us to define a table which will incorporate the overall information (seen as a set of vectors). In particular, this table represents the

Table 1 Time visiting table

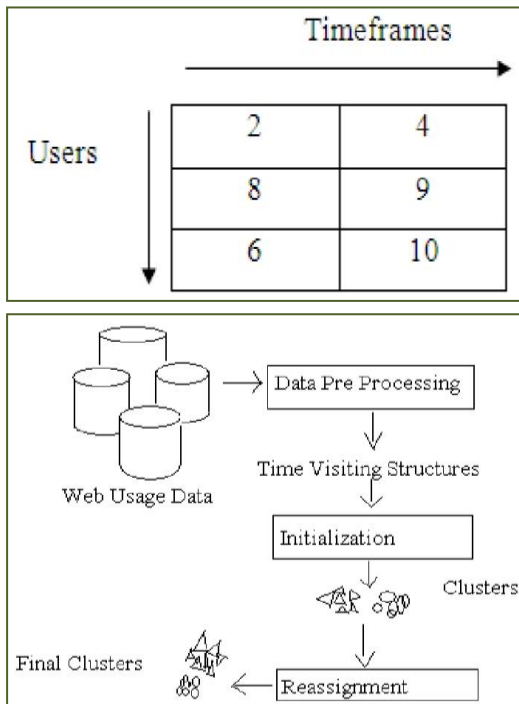


Figure 1
Clustering Process

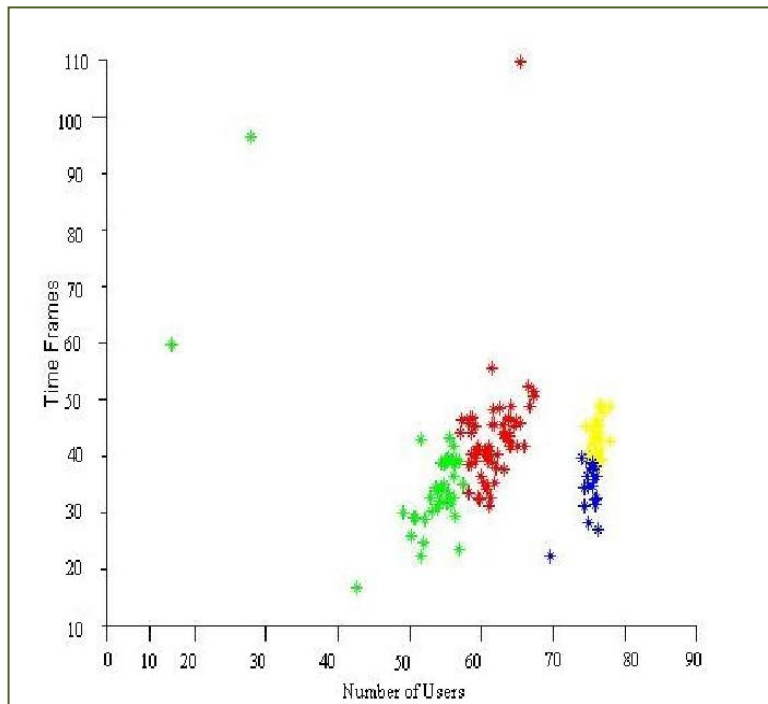


Figure 2
Result graph

frequency of visits to particular pages incorporating the exact knowledge about the timeframes of these visits. These structures are summarized next:

2.1. Time Visiting Vectors

A time visiting vector $TV(i, :)$, where $i = 1 \dots n$, represents a user's accessing behavior with respect to time (timeframes). It is also a multivariate vector consisting of t measurements:

$$TV(i, :) = (TV(i, 1) \dots TV(i, t)) \quad (1)$$

where the $TV(i, l)$ element, $l = 1, \dots, t$, indicates the number of times the user i visits the whole site (all the p pages) during the l timeframe. All the $TV(i, :)$ vectors are organized in the two dimensional $n \times t$ users' time visiting table 1. For example, in Figure 1, which depicts the table TV , the fact that $TV(2, 2) = 9$ means that the user identified as 2 has made 9 visits to the whole site during the timeframe 2.

2.2. Clustering Phases

This time related clustering algorithm is an unsupervised hard partitioned method and it is used to minimize the objective function $E_t(6)$. Normally in this clustering phase there are two steps i.e., initialization and reassignment behavior should be faced as a twofold problem that will: (i) deal with the different users' page preferences and (ii) identify the time dependencies involved in the usage navigational patterns. Thus, the problem that has to be addressed should combine the above two criteria namely the users' page preferences i.e. the page aspect and the time their visits were logged i.e. the time aspect. Since the proposed approach aims at advancing the earlier ones (which considered only the page preferences. Thus, we adopt two algorithmic approaches that differ in their initialization step and tune the two aspects based on a weight factor. The first tuning approach initiates with the page preferences and then proceeds to the time aspect while the second one follows the reverse logic.

2.2.1. Initialization

K-means partitioned clustering algorithm is used to produce the k clusters. K-means algorithm is: given n points to be clustered, a distance measure d to capture their dissimilarity and the number of clusters k to be created, the algorithm initially selects k random points as clusters' centers and assigns the rest of the $n - k$ points to the closest cluster center (according to d). Then, within each of these k clusters the cluster representative (also known as centroid or mean) is computed and the process continues iteratively with these representatives as the new clusters' centers, until convergence (Zhao et al. 2005). In this framework, given the n users and the number of k clusters are to be created. In time related algorithm the clustering considers the time aspect via time visiting structure (i.e. TV table) and uses d_t as users' dissimilarity distance measure.

2.2.2. Reassignment

The reassignment step aims at producing a CL^* clustering which enhances the initial CL to meet the two criteria of the time-aware problems. Given the initial CL clustering, it aims at finding a CL^* that minimizes the objective function E_t . More specifically, the reassignment step begins with the set of k clusters produced by the CL and involves a number of iterations. During each iteration, compute for each user u_i the fluctuation of the value of the underlying objective function (i.e. E_t) caused by moving user u_i to one of the rest $k - 1$ clusters. The reassignment phase follows K-means idea for its convergence (Srivastava et al. 2000), ending either after a number of iterations or when the objective function improvement between two

consecutive iterations is less than a minimum amount of improvement specified.

3. EXPERIMENTAL RESULT

To obtain the efficient clustering results the data set used is web log of Amazon web access data. The data set includes 98 numbers of instances, 3 attributes and 5 to 10 number of URLs. The three attributes taken are session ID, Time and Page preferences by users. The result obtained is shown in figure 2.

4. PROPOSED SYSTEM

Our defined problems are of NP-hard nature since they are a generalization of the well-known clustering problem (Garey et al. 1979) and thus we can only aim for approximate solutions. Based on the previous section, we define two algorithms to solve the TUNING TIME-AWARE CLUSTERING problem (tuning algorithms). These algorithms adopt local search heuristics which are similar in spirit with the well-known K-means algorithm (Hastie et al. 2001), which is used for our initial clustering setup. Although agglomerative does not provide approximation guarantees, it has been proved very effective in many practical problems.

4.1. Agglomerative Clustering

Hierarchical clustering algorithms are either top-down or bottom-up. Bottom-up algorithms treat each document as a singleton cluster at the outset and then successively merge (or agglomerate) pairs of clusters until all clusters have been merged into a single cluster that contains all documents. Bottom-up hierarchical clustering is therefore called hierarchical agglomerative clustering or HAC. Top-down clustering requires a method for splitting a cluster. It proceeds by splitting clusters recursively until individual documents are reached. The agglomerative clustering includes three linkage criteria namely single, complete, average linkage. The most commonly used linkage method is single linkage.

4.2. Single Linkage

In cluster analysis, single linkage, nearest neighbor or shortest distance is a method of calculating distances between clusters in hierarchical clustering. In single linkage, the distance between two clusters is computed as the distance between the two closest elements in the two clusters. The distance measure used here is Euclidean distance.

Mathematically, the linkage function – the distance $D(X, Y)$ between clusters X and Y – is described by the expression

$$D(X, Y) = \min_{x \in X, y \in Y} d(x, y) \quad (2)$$

where X and Y are any two sets of elements considered as clusters, and $d(x, y)$ denotes the distance between the two elements x and y .

Algorithm: Agglomerative Clustering

```
SIMPLE HAC( $d_1, \dots, d_N$ )
for  $n \leftarrow 1$  to  $N$ 
do for  $i \leftarrow 1$  to  $N$ 
do  $C[n][i] \leftarrow \text{SIM}(d_n, d_i)$ 
 $I[n] \leftarrow 1$  (keeps track of active clusters)
 $A \leftarrow []$  (assembles clustering as a sequence of merges)
for  $k \leftarrow 1$  to  $N-1$ 
do( $l, m$ )  $\leftarrow \arg \max \{ (i, m) : i \neq m \wedge I[i] = 1 \wedge I[m] = 1 \} C[i][m]$ 
A.APPEND( $i, m$ ) (storemerge)
for  $j \leftarrow 1$  to  $N$ 
do  $C[i][j] \leftarrow \text{SIM}(l, m, j)$ 
 $C[j][i] \leftarrow \text{SIM}(l, m, j)$ 
 $I[m] \leftarrow 0$  (deactivate cluster)
Return A
```

5. CONCLUSION

In this paper the time aware clustering approaches are introduced and evaluated: the so-called TUNING time aware clustering. The URLs in a website always have a hierarchical or tree-like structural directory. So the conversions of simple numerical features from web access are too complex.

REFERENCES

1. Qingtian Han, Xiaoyan Gao, Wenguo Wu. Study on Web Mining Algorithm Based on Usage Mining, Computer-Aided Industrial Design and Conceptual Design, 2008.CAID/CD 2008. 9th International Conference on 22-25 November, 2008
2. Yan TW, Jacobsen M, Garcia Molina H, Dayal U. From User Access Patterns to Dynamic Hypertext Linking, Proc. of Fifth International WWW Conference, 1996
3. Rajni Pamnani, Pramila Chawan Web Usage Mining: A Research Area in Web Mining, Department of computer technology, VJTI University, Mumbai
4. Garey M, Johnson D, Computers and Intractability: A guide to the Theory of NP-Completeness, New York, NY, USA: Freeman, 1979
5. Hastie T, Tibshirani R, Friedman J, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Springer, 2001
6. Bianco A, Mardente G, Mellia M, Munafo M, Muscariello L, Web user session characterization via clustering techniques, in GLOBECOM'05, IEEE, Dec. 2005, 6
7. Srivastava J, Cooley R, Deshpande M, Tan P, Web usage mining: Discovery and applications of usage patterns from web data, SIGKDD Explorations, 2000, 1(2), 12-23
8. Zhao Y, Karypis G, Topic-driven clustering for document datasets, in Proc. of the SIAM Int. Conference on Data Mining, Newport Beach, CA, USA, Apr. 2005, 358-3